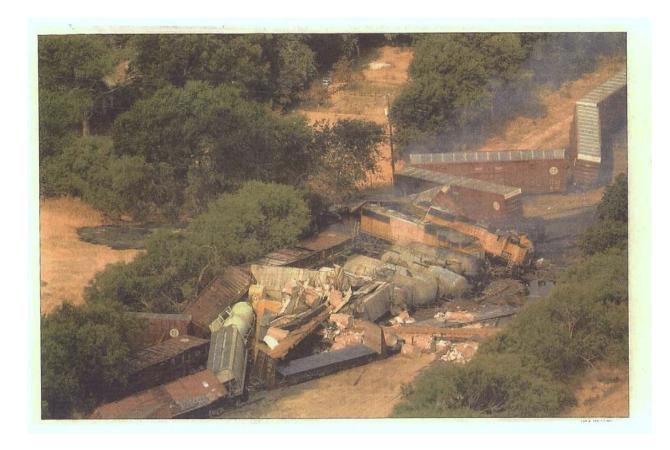


# Validation and Calibration of a Fatigue Assessment Tool for Railroad Work Schedules, Summary Report

Federal Railroad Administration

Office of Research and Development Office of Safety Washington, DC 20590



#### NOTICE

This document is disseminated under the sponsorship of the Department of Transportation in the interest of information exchange. The United States Government assumes no liability for its contents or use thereof.

#### NOTICE

The United States Government does not endorse products or manufacturers. Trade or manufacturers' names appear herein solely because they are considered essential to the objective of this report.

#### REPORT DOCUMENTATION PAGE

Form Approved OMB No. 0704-0188

Public reporting burden for this collection of information is estimated to average 1 hour per response, including the time for reviewing instructions, searching existing data sources, gathering and maintaining the data needed, and completing and reviewing the collection of information. Send comments regarding this burden estimate or any other aspect of this collection of information, including suggestions for reducing this burden, to Washington Headquarters Services, Directorate for Information Operations and Reports, 1215 Jefferson Davis Highway, Suite 1204, Arlindton, VA 22202-4302, and to the Office of Management and Buddet, Paperwork Reduction Project (0704-0188), Washington, DC 20503.

Bavis riigiway, baite 1201, riiiington, vri 22202 1002, and	to the office of Management and Baaget, 1 aper vont reducti	10111 10 Jobt (070 1 0 100), Washington, 20 20000.
1. AGENCY USE ONLY (Leave blank)	CY USE ONLY (Leave blank)  2. REPORT DATE October 31, 2006  3. REI	
4. TITLE AND SUBTITLE Validation and Calibration of a Fatigue As Summary Report	5. FUNDING NUMBERS	
6. AUTHOR(S) Steven R. Hursh, Ph.D., and Joseph F. Fanzone, Jr.	Thomas G. Raslear, Ph.D., <sup>2</sup> A. Scott Kaye, <sup>3</sup>	
7. PERFORMING ORGANIZATION NAME(S) A <sup>1</sup> Institutes for Behavior Resources <sup>2</sup> Office 2104 Maryland Avenue Feder Baltimore, MD 21218 Wash	8. PERFORMING ORGANIZATION REPORT NUMBER	
9. SPONSORING/MONITORING AGENCY NAM U.S. Department of Transportation Federal Railroad Administration Office of Research and Development Washington, DC 20590	10. SPONSORING/MONITORING AGENCY REPORT NUMBER  DOT/FRA/ORD-06/21	
11. SUPPLEMENTARY NOTES COTR: Thomas Raslear, Federal Railroad	Administration	
12a. DISTRIBUTION/AVAILABILITY STATEME	NT rough the National Technical Information Se	12b. DISTRIBUTION CODE
13. ABSTRACT (Maximum 200 words) This report summarizes the results of a	project to demonstrate a method to valid	late and calibrate a fatigue model. The

This report summarizes the results of a project to demonstrate a method to validate and calibrate a fatigue model. The project examined 30-day work histories of locomotive crews prior to 400 human factors accidents and 1000 nonhuman factors accidents. A biomathematical fatigue model estimated crew effectiveness (the inverse of fatigue) based entirely on work schedule information and opportunities to obtain sleep. A reliable linear relationship existed between crew effectiveness and the risk of a human factors accident (r = 0.93); no such relationship was found for nonhuman factors accidents. This result satisfied the criteria for model validation. A reliable time of day variation occurred in human factors accidents (r = 0.71) but not in nonhuman factors accidents. The risk of human factors accidents was elevated at any effectiveness score below 90 and increased progressively with reduced effectiveness. At an effectiveness score  $\leq 50$ , human factors accidents were 65 percent more likely than chance. Human factors accident risk increases reliably when effectiveness goes below 70, a value that is the rough equivalent of a 0.08 blood alcohol level or being awake for 21 hour following an 8-hour sleep period the previous night. Below an effectiveness score of 70, accident cause codes indicated the kinds of operator errors consistent with fatigue, confirming that the relationship between accident risk and effectiveness was meaningful

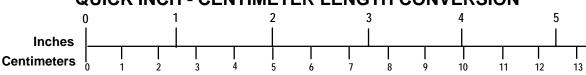
between decident risk and effectiveness was meaningful.					
14. SUBJECT TERMS	15. NUMBER OF PAGES				
Fatigue, biomathematical fat	40				
performance, effectiveness, risk, safety, freight rail operations, sleep, circadian rhythm, model validation.			16. PRICE CODE		
17. SECURITY CLASSIFICATION OF REPORT Unclassified	20. LIMITATION OF ABSTRACT				
Chelassified	Unclassified	Unclassified			

NSN 7540-01-280-5500

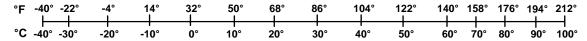
Standard Form 298 (Rev. 2-89) Prescribed by ANSI Std. 239-18 298-102

#### METRIC/ENGLISH CONVERSION FACTORS

#### **ENGLISH TO METRIC** METRIC TO ENGLISH LENGTH (APPROXIMATE) LENGTH (APPROXIMATE) 1 inch (in) = 2.5 centimeters (cm) 1 millimeter (mm) = 0.04 inch (in) 30 centimeters (cm) 1 centimeter (cm) = 0.4 inch (in) 1 foot (ft) 1 yard (yd) = 0.9 meter (m) 1 meter (m) = 3.3 feet (ft)1 mile (mi) = 1.6 kilometers (km) 1 meter (m) = 1.1 yards (yd)1 kilometer (km) = 0.6 mile (mi) **AREA** (APPROXIMATE) **AREA** (APPROXIMATE) 1 square centimeter (cm<sup>2</sup>) = 0.16 square inch (sq in, in<sup>2</sup>) 1 square inch (sq in, in<sup>2</sup>) = 6.5 square centimeters (cm<sup>2</sup>) 1 square foot (sq ft, ft<sup>2</sup>) = 0.09 square meter (m<sup>2</sup>) 1 square meter (m<sup>2</sup>) = 1.2 square yards (sq yd, yd<sup>2</sup>) 1 square kilometer (km<sup>2</sup>) = 0.4 square mile (sq mi, mi<sup>2</sup>) 1 square yard (sq yd, yd $^2$ ) = 0.8 square meter (m<sup>2</sup>) 1 square mile (sq mi, mi<sup>2</sup>) = 2.6 square kilometers (km<sup>2</sup>) 10,000 square meters $(m^2)$ = 1 hectare (ha) = 2.5 acres 1 acre = 0.4 hectare (he) = 4,000 square meters (m<sup>2</sup>) MASS - WEIGHT (APPROXIMATE) MASS - WEIGHT (APPROXIMATE) 1 gram (gm) = 0.036 ounce (oz) 1 ounce (oz) = 28 grams (gm)1 pound (lb) 0.45 kilogram (kg) 1 kilogram (kg) = 2.2 pounds (lb) 1 short ton = 2,000 pounds 1 tonne (t) = 1,000 kilograms (kg)0.9 tonne (t) = 1.1 short tons **VOLUME** (APPROXIMATE) **VOLUME** (APPROXIMATE) 1 milliliter (ml) = 0.03 fluid ounce (fl oz) 1 teaspoon (tsp) = 5 milliliters (ml) 1 tablespoon (tbsp) = 15 milliliters (ml) 1 liter (I) = 2.1 pints (pt)1 fluid ounce (fl oz) = 30 milliliters (ml) 1 liter (I) = 1.06 quarts (qt)1 cup (c) = 0.24 liter (I) 1 liter (I) = 0.26 gallon (gal) 0.47 liter (I) 1 pint (pt) = 1 quart (qt) = 0.96 liter (I) 1 gallon (gal) = 3.8 liters (l) 1 cubic foot (cu ft, ft<sup>3</sup>) = 0.03 cubic meter (m<sup>3</sup>) 1 cubic meter (m<sup>3</sup>) = 36 cubic feet (cu ft, ft<sup>3</sup>) 1 cubic meter (m<sup>3</sup>) = 1.3 cubic yards (cu yd, yd<sup>3</sup>) 1 cubic yard (cu yd, yd³) = 0.76 cubic meter (m³) TEMPERATURE (EXACT) TEMPERATURE (EXACT) $[(x-32)(5/9)] \circ F = y \circ C$ [(9/5) y + 32] °C = x °F**QUICK INCH - CENTIMETER LENGTH CONVERSION**



### **QUICK FAHRENHEIT - CELSIUS TEMPERATURE CONVERSION**



For more exact and or other conversion factors, see NIST Miscellaneous Publication 286, Units of Weights and Measures. Price \$2.50 SD Catalog No. C13 10286

## **Contents**

Illustratio	ons	V
Tables		vi
Acknowl	edgements	vii
Preface		ix
Executiv	e Summary	1
1.	Introduction	3
1.1	Background	3
1.2	Objectives	3
1.3	Overall Approach	4
2.	Method of Analysis	5
2.1	General Method	5
2.2	Validation Method	5
2.3	Sample Size Determination	7
2.4	Selection of Accidents	7
2.5	Accident Samples for Each Participating Railroad	8
2.6	Data Requested	8
2.7	The SAFTE Model and the Work Schedule Fatigue Assessment Tool	8
3.	Results of Analysis	11
3.1	Descriptive Analysis	11
3.2	Validation of a Fatigue Model	13
3.3	Calibration of a Fatigue Model	17
3.4	Other Evidence of Fatigue	
3.5	Interpretation of Effectiveness Scores	21
3.6	Interpretations and Limitations	23
4.	Summary and Conclusions	25
4.1	Summary	25
4.2	Conclusion	26
5.	References	27
Abbrevia	ations/Acronyms	29

## Illustrations

Figure 1.	Durations of Reported Job-Related Non-Sleep Intervals
Figure 2.	Frequency of Work Hours by Time of Day, Hours by Consistent Night Workers Shown as the Dark Blue Bars at the Bottom and All Others Stacked Above
Figure 3.	Accident Risk by Time of Day. Data Have Been Double-Plotted to Show the Repeating Circadian Pattern
Figure 4.	Work Interval Effectiveness Distribution for Five Railroads
Figure 5.	Proportion of Time or Accidents as a Function of Predicted Effectiveness Categories.  Human Factors and Nonhuman Factors Accidents are Indicated as Blue Triangles and Green Squares, Respectively
Figure 6.	Human Factors Accident Risk at Each Level of Effectiveness Aggregated from Five Railroads
Figure 7.	Nonhuman Factors Accident Risk at Each Level of Effectiveness Aggregated from Five Railroads
Figure 8.	Human Factors Accident Risk by Criterion Levels of Effectiveness Aggregated for Five Railroads

## **Tables**

Table 1.	Human Factors Accident Cumulative Risk at Various Criterion Levels of Effectiveness	19
Table 2.	Human Factors Accident Cause Codes Related to Effectiveness ≤ 70	20
Table 3.	The Relationship Among Various Effectiveness Scores and Other Meaningful Metric Likelihood of a Lapse, Continuous Hours Awake, and BAC	
Table 4.	The Effects of Various Daily Sleep Patterns on Effectiveness Estimates at 1600 hr an 0400 hr. Three Schedules: 1, 2, 7 Days at the Specified Sleep Level	

## **Acknowledgements**

The authors of this report wish to thank the following individuals for their invaluable contributions to the analysis, results, and documentation of the procedures:

- The Federal Railroad Administration for providing funding for this important study.
- Dr. Timothy Elsmore for the development of the batch processing software that made it possible to analyze the enormous amount of work schedule data used in this study.
- Jeff Moller and Peter French of the Association of American Railroads for encouraging this study, interfacing with industry, and assisting in the identification of a suitable sample of reportable accidents that formed the basis for the study.
- The managers and points of contacts of the five participating railroads for cooperation in conducting the study and contribution of time and effort to assemble the data needed to conduct the study.
- The Brotherhood of Locomotive Engineers and Trainmen and the United Transportation Union labor organizations for their support of the study objectives.
- Foster-Miller, Inc., the prime contractor, and Judith Gertler, the project manager, for the opportunity and generous support to conduct this study.

## **Preface**

This report is a summary of a major study sponsored by the Federal Railroad Administration (FRA) to demonstrate a method to validate and calibrate fatigue models for use in predicting and managing fatigue in railroad workers. The summary describes the methods and the major findings aggregated from data provided by five freight rail carriers. The researchers are preparing a more comprehensive report for FRA that will contain additional details of method and results. Furthermore, the researchers will make available the database used for the study for future research and analysis, with the contents coded to protect the privacy of the participants.

## **Executive Summary**

Biomathematical fatigue models allow the objective assessment of fatigue so that employees and employers can schedule work and rest to minimize the degradation of operator performance by fatigue. To be useful, a fatigue model must be validated. Validation means that the model must be a predictor of fatigue-related performance errors. Moreover, a model should be calibrated. Calibration means that the predictions from a model can be related to the level of risk of failures of human performance. One method of validating and calibrating a biomathematical fatigue model is to demonstrate that the model can predict an increased likelihood of human factors accidents relative to nonhuman factors accidents under conditions of fatigue. A valid fatigue model should predict higher levels of fatigue (based on opportunities to sleep and an accident's time of day) when a greater likelihood of a human factors accident exists. By comparison, fatigue levels should have a weaker or no relationship to the likelihood of nonhuman factors accidents. The Federal Railroad Administration (FRA) Office of Research and Development and the Office of Safety have partnered with the railroad industry to demonstrate a method to validate and calibrate fatigue models. This study collected 30-day work histories of locomotive crews prior to 400 human factors and 1000 nonhuman factors accidents to demonstrate this validation method. A total of over 1 million 30-minute work intervals before the accidents, covering over 57,000 work starts, were evaluated for effectiveness (the inverse of fatigue) predicted by the Sleep, Activity, Fatigue, and Task Effectiveness (SAFTE) model using the Fatigue Avoidance Scheduling Tool (FAST). This served as the baseline of exposure to various levels of effectiveness. In addition, the effectiveness at the time of each accident was calculated.

The analysis used two criteria to evaluate if the SAFTE biomathematical fatigue model was a valid predictor of fatigue-related accidents:

- The proportion of human factors accidents that occur at low levels of effectiveness should be greater than the proportion of time working at those levels of effectiveness (exposure level). If the proportion of human factors accidents is reliably greater than the exposure level, and there is a statistically reliable relationship between decreases in effectiveness and an increase in human factors accident risk, then low levels of effectiveness (increased fatigue) predict increased human factors accident risk.
- By comparison, there should not be a statistically reliable relationship between
  effectiveness and nonhuman factors accident risk, and at low levels of effectiveness
  (increased fatigue), the risk of human factors accidents should be greater than the risk of
  nonhuman factors accidents. If the human factors accident risk is reliably greater than
  the nonhuman factors accident risk, then low effectiveness is associated with the kinds of
  accidents that would be expected to be related to fatigue.

The statistical reliability of relationships was based on finding significant correlation coefficients (r). The chi square ( $\chi^2$ ) statistic was used to compare the distributions of human factors and nonhuman factors accidents to demonstrate significantly different distributions of risk. In each case, the researchers rejected the null hypothesis (the hypothesis that there is no correlation (i.e., r = 0) or that there is no difference in the distributions of human factors and nonhuman factors accidents) when the chance probability of a finding was less than 5% (p < 0.05).

The results of the study indicated that the biomathematical fatigue model met both criteria for validation and the results were statistically reliable:

- The biomathematical fatigue model estimates of effectiveness were significantly correlated with human factors accident risk (r = -0.93, p < 0.01, Figure 6).
- The biomathematical fatigue model distinguished between human factors and nonhuman factors accidents. Nonhuman factors accident risk was not correlated with estimated effectiveness (r = -0.14, p > 0.05, Figure 7), and at low estimated effectiveness the relative risk of human factors and nonhuman factors accidents were significantly different ( $\chi^2 = 7.201$ , p < 0.01).

This study found that an elevated risk of human factors accidents occurs at any effectiveness score below 90, and accident relative risk increased as effectiveness decreased. A reliable increase in human factors accident risk occurred when effectiveness scores were below 70 but nonhuman factors accident risk was not consistently elevated. Based on other research, effectiveness scores below 70 are the rough equivalent of a 0.08 blood alcohol level or being awake for 21 hours (hr) following an 8-hr sleep period the previous night.

An analysis of the cause codes associated with accidents that occurred at or below an effectiveness score of 70 showed an overrepresentation of the type of human factors accidents that might be expected of a fatigued crew (e.g., failure to comply with a stop signal). This confirms that the relationship between reduced effectiveness and elevated human factors accident risk is meaningful and not a mere correlation.

The study also demonstrated that human factors accidents follow a circadian pattern that is significantly correlated with the circadian rhythm of a fatigue model (r = 0.71, p < 0.05). The same model rhythm is not correlated with nonhuman factors accidents. The maximum human factors accident risk due to time of day alone was increased by less than 20 percent, while the maximum accident risk due to reduced effectiveness (fatigue) was increased by 65 percent, reflecting the combined effects of time of day and sleep opportunities.

This study was designed to demonstrate a method to test the validity of fatigue models. The data used accidents and the 30-day work histories that preceded them were not a random sample of all workers in freight rail operation. Hence, the levels of effectiveness calculated by the SAFTE model should not be interpreted as representative of the freight railroad work force in general. The study was not designed to determine the extent of fatigue in the freight rail industry. Furthermore, given the well-known variations in individual sleep requirements and absence of specific information on individual sleep habits, health, and circumstances, it was not the intent of this study to validate fatigue models based entirely on work schedule data as tools for determining the fatigue of particular individuals.

This study provides the first evidence that a biomathematical fatigue model can relate work schedules to an elevated risk of railroad accidents. This provides a strong scientific basis for evaluating work schedules with valid mathematical models to reduce worker fatigue. A mathematical model for detecting elevated fatigue risk could be part of a nonprescriptive, performance-based fatigue management plan that would supplement current regulations. Although fatigue models do not identify all sources of fatigue and will require a cooperative partnership among management, labor, and government regulators, they are an important tool in the identification of one of the causes of fatigue in the railroad industry.

#### 1. Introduction

FRA has sponsored a research program to demonstrate a method to validate and calibrate fatigue models for use in predicting and managing fatigue in railroad workers. Fatigue models allow the objective assessment and forecasting of fatigue so that employees and employers can schedule work and rest to minimize degradation of operator performance by fatigue. To be useful, a fatigue model must be validated. Validation means that a fatigue model predicts changes in job performance and/or job-related errors, such as incidents and accidents, caused by fatigue. A useful fatigue model must be calibrated. Calibration means that the predictions from the model can be related to the risks of meaningful failures of human performance. This report describes the results of the research program to demonstrate a method to validate and calibrate fatigue models.

#### 1.1 Background

Human factors accidents have increased as a proportion of FRA reportable accidents over the past 5 years. No question exists that fatigue may be a factor in many human factors accidents. Without a detailed history of work and rest before an accident, however, it is difficult to determine the role of fatigue in that accident. The analysis of work/rest histories to rule out fatigue can be accomplished with a number of software models associated with fatigue, but until now none of these models has been validated in commercial transportation operations. In the interest of developing validated fatigue models, the Office of Research and Development and the Office of Safety sponsored a study of work histories of locomotive crews associated with accidents to provide the necessary data to validate fatigue models. This report describes the details of the data collection, the methods used to analyze the data for signs of fatigue, and the results of the analysis.

#### 1.2 Objectives

One method of validating and calibrating a fatigue model is to demonstrate that the model can predict an increased likelihood of human factors accidents relative to nonhuman factors accidents under conditions of fatigue. A valid fatigue model should predict higher levels of fatigue (based on opportunities to sleep and the time of day of an accident) when a greater likelihood of a human factors accident exists. By comparison, fatigue levels should have a weaker or no relationship to the likelihood of nonhuman factors accidents. The present study determined cognitive effectiveness (a predictor of speed of reactions and vigilance in laboratory tests that is inversely related to fatigue) from 30-day work histories of locomotive crews prior to 400 human factors and 1000 nonhuman factors accidents. The objective was to determine if a statistically reliable relationship exists between reductions in effectiveness and the risk of human factors accidents. Further, the study determined if the relationship of effectiveness to human factors accidents was larger and more consistent than that of nonhuman factors accidents, which would be expected to be much less sensitive to the effects of fatigue. The second objective was to determine the level of effectiveness at which an elevated likelihood of human factor accidents occurs relative to chance. That result served to calibrate the fatigue model for aggregate work schedule analyses of fatigue. The goal was to determine the nature of the relationship between effectiveness scores and statistically elevated accident risk. Given the well-known variations in individual sleep requirements and absence of specific information on individual sleep habits,

health, and circumstances, it was not the intent of this study to validate fatigue models based entirely on work schedule data as tools for determining the fatigue of particular individuals.

The biomathematical fatigue model used for these studies, the SAFTE model, was originally developed for the U.S. Army and Air Force to predict potential fatigue from work schedules (Hursh, Redmond, Johnson, Thorne, Belenky, Balkin, Storm, Miller, and Eddy, 2004). The model has been incorporated into a fatigue assessment tool (Hursh, Balkin, Miller, and Eddy, 2004) called FAST. FAST uses work schedule information to estimate the amount of sleep and cognitive effectiveness. Cognitive effectiveness is a metric that tracks speed of performance on a simple reaction time test and is strongly related to overall response speed, vigilance, and the probability of lapses (Hursh, et al., 2004; Van Dongen, 2004). Cognitive effectiveness can be interpreted as the inverse of fatigue.

#### 1.3 Overall Approach

FRA and the Association of American Railroads (AAR) determined which human factors accident cause codes would be appropriate to use. This was important since work histories for locomotive crews would be analyzed to determine effectiveness, and, therefore, only cause codes that related to locomotive crew errors would be relevant. FRA and AAR also determined which nonhuman factors accident cause codes would be appropriate to use as a control group. Finally, the kinds of track and equipment most likely to be associated with irregularly scheduled line haul (between cities or terminals) freight service was determined. For each accident, the researchers requested from that railroad the 30-day work history of the train crew just before the accident. The work histories were used to estimate the effectiveness of the crew at the time of the event and to determine the overall percentage of time typically spent at work at various levels of effectiveness. This information was used to determine the ability of the fatigue model to predict increases in human factors accident risk relative to chance. It was also used to determine if the relationship between effectiveness and accident risk was stronger for human factors than nonhuman factors accidents. Little or no relationship should exist between effectiveness and nonhuman factors accident risk. The reduction in effectiveness associated with statistically reliable cumulative increases in accident risk was determined in order to calibrate the effectiveness scores. Additionally, human factors accidents that were overrepresented at reduced levels of effectiveness were examined to determine if these accidents were consistent with the expected effects of crew fatigue (e.g., lapses increase with fatigue and are consistent with accidents in which signals are passed at danger). Descriptive information regarding work patterns and accident patterns that were independent of the model predictions were also developed.

## 2. Method of Analysis

#### 2.1 General Method

The general method of the study was to use a biomathematical fatigue model (SAFTE) to estimate the cognitive effectiveness of locomotive crews at the time of accidents to determine if low levels of effectiveness are associated with higher than expected rates of accidents. The model computes cognitive effectiveness (the inverse of fatigue) from opportunities to sleep and the time of day during the 30 days before the accident and at the time of the accident. To control for extraneous factors, the analysis compared human factors and nonhuman factors accidents using the same methods. The following sections describe the validation method, the sample determination, the data requested from the railroads that describe the accidents, and the 30-day work histories of the associated train crews.

#### 2.2 Validation Method

Thirty-day work histories were collected for train crews who were involved in FRA reportable accidents. FAST processed these histories with the SAFTE model to determine the crew's cognitive effectiveness at the time of the accidents. The histories were also used to estimate the proportion of time spent working at various levels of effectiveness during the preceding 30 days. This is the exposure of the locomotive crews to fatigue. If fatigue is not associated with accidents, then the proportion of accidents at each effectiveness level should be similar to fatigue exposure. Fatigue exposure provides an expected proportion of accidents at each effectiveness level. An easy way to see if the obtained proportion of accidents matches the expected proportion is to divide the obtained proportion by the expected proportion. If accidents are independent of fatigue, the ratio should be close to one. This ratio is the accident risk. Statistically reliable deviations in the ratio indicate increased or decreased risk, depending on whether the ratio is greater or less than one. This analysis yields a function that relates accident risk to levels of effectiveness and estimates fatigue-associated risk.

It was expected that if fatigue is a factor in some accidents, then accident risk should be higher at lower levels of effectiveness (increased levels of fatigue), and the risk should be greater for human factors accidents than for nonhuman factors accidents. Testing these hypotheses takes several steps and makes several assumptions.

#### 2.2.1 Assumptions

First, it is assumed that loss of alertness, confusion, or slowed reaction time causes some railroad accidents. These are all performance changes that are associated with fatigue. Since these accidents are the result of human error, it was expected that fatigue-associated accident risk would be greater for human factors than for nonhuman factors accidents.

Second, it is assumed that when the model predicts a loss in cognitive effectiveness, an increase occurs in the probability of a human error that could result in a railroad accident. For example, when the model predicts reduced effectiveness, the train crewmen are more likely to miss a signal or run through a switch. This connection is not deterministic; reduced effectiveness increases the probability of an error but does not determine that an error will occur. Furthermore, a similar error could occur in the absence of unusual fatigue, such as caused by a distraction, so that not all occurrences of a particular sort of error can be attributable to fatigue.

Hence, a fatigue model can predict an increase in fatigue-associated risk, not the specific occurrence of an accident.

Third, the SAFTE model incorporated in FAST makes a number of assumptions about the amount of sleep that can be achieved under specific work schedules, the amount of sleep the average person requires, and the susceptibility of the average person to time of day variations in alertness. Wide variations exist in individual sensitivity to the factors that cause fatigue, so, again, fatigue models can only predict an increased risk of fatigue, not a specific individual person's level of fatigue or performance.

#### 2.2.2 Validation Criteria

If the fatigue model is valid, it will show that human factors related accidents are more likely to occur at decreased levels of crew effectiveness, while the likelihood of nonhuman factors-related accidents are the same regardless of crew effectiveness.

The fatigue model estimates effectiveness during each work shift of each worker in the database up to the time of an accident and records the mean effectiveness level across each 30-minute work interval. From all work histories combined, the approximate proportion of work time spent by workers at given effectiveness levels can be computed. This is the exposure to fatigue and, as noted above, determines the expected proportion of accidents at any given effectiveness level. The model also estimates the crew effectiveness at the time of the accident. From these data the proportion of accidents occurring at any given crew effectiveness level can be computed. Combining the two proportions defines the risk ratio:

$$Risk\ Ratio = \frac{\left(Accidents\ at\ Effectiveness\ Level\right)/\left(Total\ Number\ of\ Accidents\right)}{\left(Work\ Time\ at\ Effectiveness\ Level\right)/\left(Total\ Work\ Time\right)}$$

Accident risk, as used in this report, is defined entirely in terms of the risk ratio specified above. If operator effectiveness is unrelated to accident occurrence, then the proportion of accidents occurring at any given level of effectiveness should be approximately equal to the proportion of work time spent at that effectiveness level (i.e., the accident risk should be close to 1). If lower effectiveness is related to a greater likelihood of an accident, then accident risk should increase with decreasing effectiveness (i.e., be greater than 1 at lower effectiveness levels and less than 1 at higher effectiveness levels).

Therefore, model-derived estimates of effectiveness that show the following would constitute validation that the model provides a measure of fatigue-associated accident risk:

- 1. A statistically reliable relationship between decreases in effectiveness and an increase in human factors accident risk
- 2. The absence of a statistically reliable relationship between effectiveness and <u>nonhuman</u> factors accident risk, and a greater risk of human factors versus nonhuman factors accidents at low levels of effectiveness

The statistical reliability of relationships was based on finding significant correlation coefficients (r). The chi square ( $\chi^2$ ) statistic was used to compare the distributions of human factors and nonhuman factors accidents to demonstrate significantly different distributions of risk. In each case, we rejected the null hypothesis (the hypothesis that there is no correlation, i.e., r = 0, or that there is no difference in the distributions of human factors and nonhuman factors accidents) when the chance probability of a finding was less than 5 percent (p < 0.05).

#### 2.2.3 Calibration of the Fatigue Model

Up to this point, the analysis method has referred to low levels of effectiveness and levels of fatigue without specifying what those levels might be. This is because it is not known in advance how much fatigue or reductions in cognitive effectiveness are sufficient to cause a detectable increase in accidents or a detectable difference between human and nonhuman factors accidents. The level of effectiveness for human factor accidents that yields a reliable increase in accidents relative to chance and the effectiveness level for human factor accidents that is reliably different from the level for nonhuman factors accidents were used to calibrate the model.

#### 2.2.4 Variables Not Considered

Given the nature of this analysis of historical work times associated with reportable accidents. the biomathematical fatigue model was not able to consider some variables that might ordinarily be part of an analysis. The database included work start and stop times, typical call times for each location, and estimated average commute times for each location. The analysis could not consider each worker's usual sleep habits or actual sleep times, quality of the sleep environment, or schedule predictability. The analysis did not consider the use of fatigue countermeasures or additional naps. The analysis could not consider unusual events that might have interfered with usual sleep opportunities. The analysis could not consider schedule delays or misinformation, medical conditions and sleep disorders, medications and/or drug use, observations of operator performance and appearance, concurrent stress, family issues, work demands, or crew resource management issues, such as communication problems. Many of these factors could cause additional fatigue or performance disruptions. Some of these factors could reduce fatigue, such as the use of napping strategies. The analysis had to assume that each individual obtained as much sleep as was afforded by the schedule, call times, commute times, and personal obligations, up to a nominal level of 8 hrs of sleep per day. Given these limitations, the predicted levels of effectiveness reported in this study are not necessarily the same as might be obtained if an accident investigator performed a detailed fatigue analysis at the time of the accident and caution should be exercised in extending the calibration levels to that application.

#### 2.3 Sample Size Determination

A statistical process called power analysis, based on the results of a pilot study to evaluate the odds of finding a statistically significant difference in fatigue between human factor and nonhuman factor accidents given various sample sizes and assumptions about the signal to noise ratio, was used to determine sample size. The power analysis indicated that sufficient accidents occur within a 2½-year period to reach an acceptable statistical power with a sample size of 400 human factors cases and 1000 nonhuman factors cases. Since usually two employees are on each locomotive, an engineer and conductor, this results in approximately 2800 work histories and accident cases.

#### 2.4 Selection of Accidents

Accidents that were assigned a human factors cause code (Hxxx) for either the primary or secondary cause were considered human factors accidents. Human factors accidents caused by actions of non-operational personnel, such as maintenance or signal workers, cannot logically be associated with the fatigue or effectiveness levels of the crew operating the train. Accidents with cause codes of H305, H402, H501, H993, H994, and H997 were, therefore, excluded from the

list of human factors accidents considered. All track- and equipment-related accident cause codes (Txxx or Exxx) were included for the nonhuman factors accidents.

#### 2.5 Accident Samples for Each Participating Railroad

To achieve the sample size specified in Section 2.3, it was estimated that data would be needed for all non-excluded human factors accidents occurring during calendar years 2003, 2004, and the first 6 months of 2005 on mainline, siding, and industry track involving freight trains. All track- and equipment-caused accidents occurring over the same period on mainline, siding, and industry track involving freight trains were used as a comparison set of nonhuman factors accidents. The 5 participating railroads—BNSF, CSX, Kansas City Southern (KCS), Norfolk Southern (NS), and Union Pacific (UP)—provided data on a total of 405 non-excluded human factors accidents and 1015 nonhuman factors accidents reported in this interval, which constituted the study sample.

#### 2.6 Data Requested

Each participating railroad provided work schedule data for each worker for the 30 days before each of the identified accidents. A template file was provided for the railroads to use to construct the data tables and is described in the full report. Most critical to the analysis was a record of the reporting and release date/times of each work episode, including deadhead and limbo times, which for analytical purposes were combined as unavailable for sleep. (The duty periods reported are not all regulated by hours of service rules and may exceed the allowable limit of 12 hr on duty.) Data on terminals, standard call times, and estimated commute times were also requested and applied to further limit opportunities to sleep. Worker identification was used solely to associate each work history with the accident involved and was coded to insure confidentiality.

### 2.7 The SAFTE Model and the Work Schedule Fatigue Assessment Tool

The biomathematical fatigue model used for this specific analysis was the SAFTE model. The U.S. Air Force and FRA have sponsored the development of a scheduling tool based upon the model that can be used to assess and manage fatigue in aircrews and railroad workers. The software is called FAST. For assessment of aggregate fatigue in a work force, FRA sponsored the development of a work schedule fatigue assessment tool based on FAST that can process hundreds of work histories compiled in a database and provide aggregate estimates of fatigue across workers and for specific times associated with accidents. This project used the batch processing version of FAST.

## 2.7.1 The Work Schedule Fatigue Assessment Process

In general, the batch processing involves the following steps. First, the data are entered into a set of databases. Second, the program processed those work histories to determine when sleep could have occurred, given that workers normally cannot sleep at work, generally do not sleep after they are called to work, and must commute to and from their place of work. The tool uses an algorithm called AutoSleep to estimate how much sleep the average railroad worker would get given the sleep opportunities afforded by the work schedule and the call/commute times. AutoSleep is based on a study of railroad engineer work/rest diaries (Pollard, 1996) and is described in more detail in the full report. The third step was to estimate the cognitive effectiveness of each worker for every minute that

the worker was awake, using the SAFTE model. A detailed description of the SAFTE model has been previously published (Hursh, et al., 2004). Simply stated, crew effectiveness is an estimate of the speed of reaction time and alertness. Estimated effectiveness varies with the combined effect of the time of day and the pattern of sleep. Crew effectiveness follows a daily (circadian) rhythm that is much lower between 0000 and 0600 hr than between 1200 and 1800 hr. In addition, the model keeps track of the amounts of sleep opportunities a person gets as it contributes to performance and progressively degrades performance if the person experiences a loss of sleep from the nominal requirement of 8 hr per day. The batch processor summarized these detailed effectiveness estimates into 30-minute averages during every work period in the database. The program also summarized effectiveness in 30-minute intervals around the clock to provide a time of day estimate of work effectiveness. The analysis summarized effectiveness for each person's work shift by describing the distribution of the amount of time spent at work as a function of effectiveness. Finally, and most importantly, the program estimated the effectiveness of each worker at the time he/she was reported to have been involved in an accident.

#### 2.7.2 The Dimensions of the Analysis

The batch processor created a schedule file for each worker that could be viewed in the standard version of FAST. These schedule files were used to troubleshoot any unusual values and in several cases led to corrections to the database. Of the 2962 workers involved in accidents (approximately 2 crewmembers per accident) that were reported, valid work histories were available for 2843 workers, a 96 percent success rate. Of those, 790 were involved in human factors accidents, and 2053 were involved in nonhuman factors accidents. Each work history was composed of a series of records that constituted a work shift. Across the five railroads, the analysis processed 57,537 work shifts. In total, the results reported are based on over 1 million 30-minute work intervals. Hence, the effectiveness exposure estimates and accident effectiveness values were based on over 2800 work histories, and the results of all crewmembers contributed equally to the findings. The data were not coded in such a way to permit estimates of effectiveness by craft (engineer or conductor).

-

<sup>&</sup>lt;sup>1</sup> Crew effectiveness can have values in the range of 0 to 100. A person who consistently obtains 8 hr of good quality sleep would have a peak effectiveness score of 100 during the following waking period.

### 3. Results of Analysis

### 3.1 Descriptive Analysis

The first results are simple descriptions of various features of the work schedule and accident data without any reference to the fatigue model predictions of effectiveness. The three basic descriptive charts are:

- Work Durations
- Clock Time of Work
- Accidents by Time of Day

For this and all other sections of the report, the results are shown as an aggregate of the findings based on the data from all railroads.

#### 3.1.1 Job-Related Non-Sleep Durations

The railroads reported all times crewmembers were doing work related activities, such as performing as a crewman, deadheading, or in limbo time before release. Together, these constituted job-related times unavailable for sleep. Figure 1 shows the number of these intervals sorted by duration in hourly intervals. Relatively few intervals were shorter than 5-hr long. Work periods were about equally distributed from lengths of greater than 6 hr to less than 11 hr.

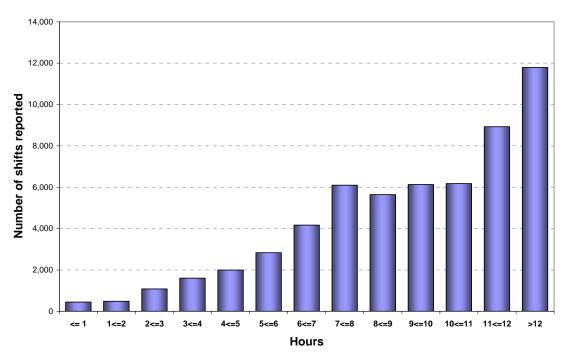


Figure 1. Durations of Reported Job-Related Non-Sleep Intervals

The most frequent job-related intervals were those between 11 and 12 hr and longer than 12 hr. This may not be surprising considering that these intervals would include deadhead time and limbo time. For example, if a typical deadhead time is at least 3 hr

long, then duty shifts of 9 hr or longer that were combined with deadhead time would all fall in the longer than 12 hr category.

#### 3.1.2 Clock Time of Work

Figure 2 shows the number of work hours that occurred at hourly intervals around the clock. This analysis indicates the frequency of work at varying times of day and shows that, in general, work was fairly evenly distributed around the clock. For the population at large (the full height of the bars in Figure 2), work is a bit more concentrated between 0800 and 1600 hr and diminishes between 1600 and 0500 hr. A detailed analysis revealed that in the overall population involved in accidents, a subgroup of workers existed who started a majority of their work periods between 2200 and 0400 hr. These workers are termed consistent night workers. Their data were separated from the larger population and evaluated as a distinctly different group. Their distribution of work is understandably more concentrated between 0000 and 1200 hr, shown as the dark blue bars in the graph. The remainder of this report focuses on the findings with those who were not consistent night workers. The full report includes the analysis of night workers.

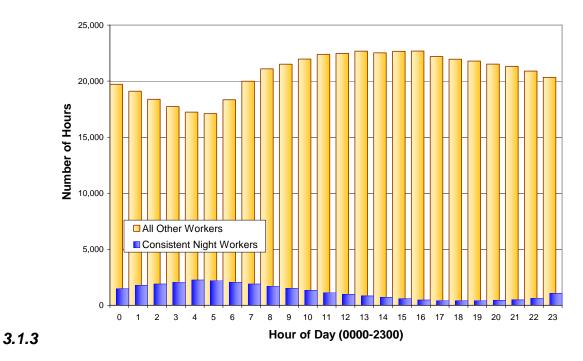


Figure 2. Frequency of Work Hours by Time of Day, Hours by Consistent Night Workers Shown as the Dark Blue Bars at the Bottom and All Others Stacked Above

### 3.1.4 Accidents by Time of Day

Figure 3 shows the risk of human factors and nonhuman factors accidents in 3-hr intervals around the clock. For each clock interval the analysis calculated the proportion of accidents that occurred in that interval and divided that proportion by the proportion of work times that occurred at that clock time based on the work interval values in Figure 2. If accidents were randomly distributed around the clock, then the points would have a value of 1.0 and fall on the dashed line. The blue triangles in Figure 3 are for human factors accidents, and the green squares are for nonhuman factors accidents. A clear

circadian pattern of human factors accident risk exists. Taking into consideration the distribution of work around the clock, accidents are relatively more likely in the early morning hours from 0000 to 0300 hr (the circadian nadir) and in the early afternoon from 1200 to 1500 hr (the postprandial dip). Accidents are much less likely in the late morning (0900 to 1200 hr) and in the early evening (1500 to 1800 hr). These patterns are predicted by a circadian pattern generated by the SAFTE model, shown as a red line. Since effectiveness is thought to decline with increased fatigue, the inverse of effectiveness is plotted against the right-hand axis. The maximum value of inverse effectiveness corresponds to the early morning peak in accident risk. The data are plotted twice along the x-axis to illustrate the rhythmic pattern of the results. Human factors time of day accident risk is reliably correlated with the circadian pattern derived from the fatigue model (r = 0.71, p < 0.05). Circadian rhythms account for 51 percent of the variance in the time of day at which human factors accidents occur. By comparison, circadian rhythms only account for 6 percent of the variance in the time of day at which nonhuman factor accidents occur.

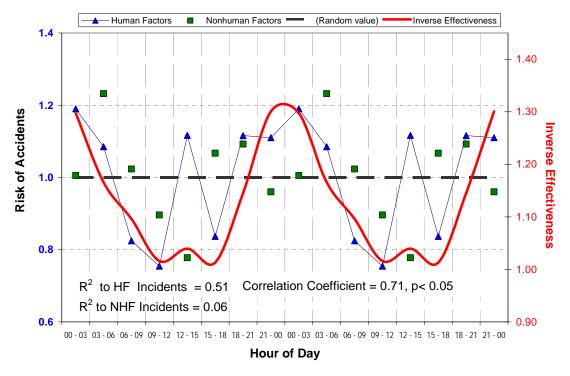


Figure 3. Accident Risk by Time of Day. Data Have Been Double-Plotted to Show the Repeating Circadian Pattern

#### 3.2 Validation of a Fatigue Model

## 3.2.1 Work Time in Effectiveness Categories

Figure 4 shows the distribution of 30-minute work intervals as a function of effectiveness. The figure shows the proportion of time spent at work with an effectiveness score between 90 and 100, between 80 and 90, between 80 and 70, and so on down to the lowest bin with scores of 50 or less effectiveness. The right-hand axis shows that 42 percent of the time locomotive crews have effectiveness scores above 90.

Another 23 percent of the time, they have effectiveness scores between 80 and 90. Hence, 65 percent of the time, locomotive crews have scores above 80. The percent of time declines consistently with decreasing effectiveness, and less than 5 percent of the time effectiveness is below a score of 50.

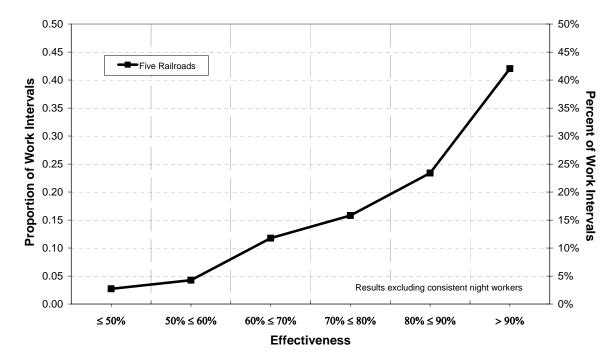


Figure 4. Work Interval Effectiveness Distribution for Five Railroads

These data provide an important reference for interpreting the estimated effectiveness values associated with accidents. If accidents are independent of effectiveness (or fatigue), then when accidents are sorted by crew effectiveness at their time of occurrence, the accidents should distribute exactly as work time effectiveness is distributed in Figure 4. On the other hand, if accidents are caused, in part, by low effectiveness (fatigue or lack of alertness), then one would expect that a greater proportion of human factors accidents would occur at low levels of effectiveness than the proportion of time or exposure to those levels of effectiveness. In other words, if the fatigue model is predictive of accidents, then low effectiveness should be associated with an elevated risk of an accident.

#### 3.2.2 Effectiveness at the Time of the Accident

Figure 5 shows the distributions of human factors and nonhuman factors accidents by effectiveness scores of the locomotive crews at the time of the accidents. Figure 5 also shows the mean work interval effectiveness distribution from Figure 4 for comparison. The y-axis is logarithmic, so that the distance between the accident distribution lines and the time distribution line at any effectiveness level reflects the degree of risk at that level; the greater the distance between the lines, the larger the ratio that defines risk (see Section 2.2.2).

Figure 5 shows that the proportion of accidents that occur above an effectiveness score of 90 is less than expected by the distribution of time, indicating reduced risk. Below an effectiveness score of 90, the rate of human factors accidents is consistently above the work-time distribution (heavy line) and gradually separates from the line as effectiveness decreases. For nonhuman factors accidents, the proportion is sometimes above the line and sometimes below the line, and no consistent relationship exists between nonhuman factors proportions and decreases in effectiveness.

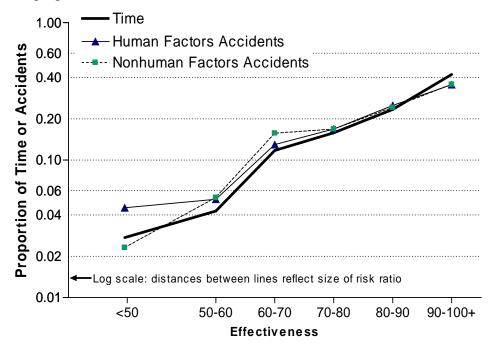


Figure 5. Proportion of Time or Accidents as a Function of Predicted Effectiveness Categories. Human Factors and Nonhuman Factors Accidents are Indicated as Blue Triangles and Green Squares, Respectively

#### 3.2.3 Accident Risk as a Function of Effectiveness

As discussed in Section 2.2.2, if effectiveness (fatigue) is not associated with accidents, then the ratio that defines relative risk will be equal to 1.0. In other words, if being exposed to a certain level of effectiveness (fatigue) does not alter chances of having an accident, then accident risk is approximately 1.0. On the other hand, a ratio greater than 1.0 indicates that accidents are more likely than chance at that level of effectiveness. A ratio of 1.5 means that a 50 percent increase in the risk of having an accident occurred at that level of effectiveness. To show the relationship between risk and effectiveness, the data points from Figure 5 were used to compute the risk ratio (using the expression from Section 2.2.2) for each effectiveness category. Figures 6 and 7 show the accident risk for the aggregated data from all five railroads for human factors and nonhuman factors accidents, respectively.

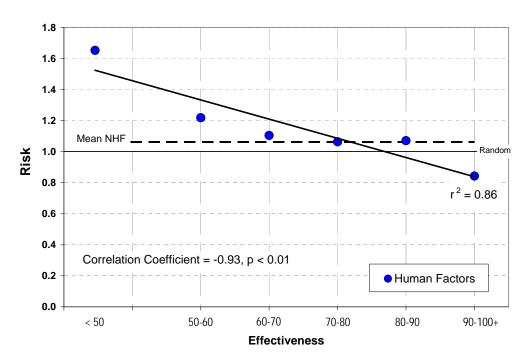


Figure 6. Human Factors Accident Risk at Each Level of Effectiveness Aggregated from Five Railroads

In Figure 6, the solid blue circles and the heavy solid line fit to them show the increase in risk for human factors accidents as a function of decreases in effectiveness (increasing fatigue) at the time of the accident. A risk value of 1.0 indicates no effect of fatigue on accident occurrence. The dashed line shows the mean accident risk of nonhuman factors accidents. As can be seen from Figure 6, human factors accident risk is well described by a linear function. The line that is fit to these data accounts for 86 percent of the variance in human factors accident risk, and a significant inverse correlation exists between accident risk and effectiveness, the inverse of fatigue (r = -0.93, p < 0.01). The data show that for effectiveness scores between 90 and 100 (values associated with optimal prior sleep), a reduction of risk occurs. At effectiveness scores below 90, risk progressively increases. At the lowest level of effectiveness, a 65 percent increase in accident risk occurs (for effectiveness scores equal to or less than 50). The consistent relationship between reduced effectiveness (increased fatigue) and elevated risk indicates that the additional risk is associated with fatigue.

• This finding satisfies the first criterion for model validation: There was a significant correlation between model predicted reductions in effectiveness and an increase in human factors accident risk.

Importantly, the maximum increase in accident risk due to time of day alone was less than 20 percent (Figure 3), while the maximum increase in accident risk due to reduced effectiveness (fatigue) was 65 percent (Figure 6), reflecting the combined effects of time of day and sleep opportunities.

Figure 7 shows the risk for nonhuman factors accidents as a function of effectiveness. In Figure 7 no consistent, statistically reliable relationship exists between risk and effectiveness. Data points fall above and below the line, showing a risk value of 1.0.

The best fitting straight line has a slope that is not reliably different from zero, indicating that no consistent relationship exists between nonhuman factors accident risk and effectiveness (r = -0.14, p > 0.05). Only 2 percent of the variance in risk is accounted for by a linear function.

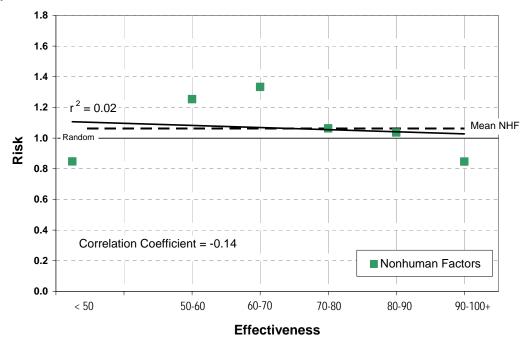


Figure 7. Nonhuman Factors Accident Risk at Each Level of Effectiveness Aggregated from Five Railroads

A large difference in relative risk occurred between human factors and nonhuman factors accidents in the lowest category of effectiveness, and this difference was statistically reliable ( $\chi^2 = 7.201$ , p < 0.01). The full report includes details of the statistical analysis.

• These findings satisfy the second criterion for model validation: There was no reliable relationship between effectiveness and nonhuman factors accident risk, and at low levels of effectiveness, human factors accident risk was elevated more than nonhuman factors accident risk.

#### 3.3 Calibration of a Fatigue Model

#### 3.3.1 Accident Likelihoods at Selected Criterion Levels

In addition to validating that a fatigue model can associate reductions in effectiveness with a reliable increase in human factors accidents, this study sought to calibrate the level of reduced effectiveness below which there is a reliable increase in accident risk (cumulative risk). To do this, a somewhat different analysis was conducted. In this case, a criterion level was set to include the proportion of all accidents that occurred with an effectiveness score at or below that level. The analysis then compared that proportion to the proportion of exposure time that occurred at or below that level. If the ratio of the two values is reliably greater that 1.0 at a particular criterion level, then that level could be considered the effectiveness level below which an increased cumulative risk of

accidents would exist. As before, the analysis was conducted on the aggregated data, excluding the consistent night workers. Figure 8 shows the results, while Table 1 summarizes them. The blue symbols are the cumulative human factors fatigue-associated accident risk, the black line is the expected cumulative risk of accidents if they are distributed as work interval effectiveness exposure (see Figure 4), and the dashed line is the mean cumulative risk ratio for nonhuman factors accidents.

Each point shows the 95 percent confidence limits. If the lower confidence limit is above the black line (risk = 1), then the bar is reliably greater than 1.0 (chance). Interestingly, for any effectiveness criterion score below 90, human factors accidents have a cumulative risk reliably greater than expected by chance. Below an effectiveness score of 70, human factors cumulative risk is reliably greater than nonhuman factors cumulative risk and reliably greater than chance.

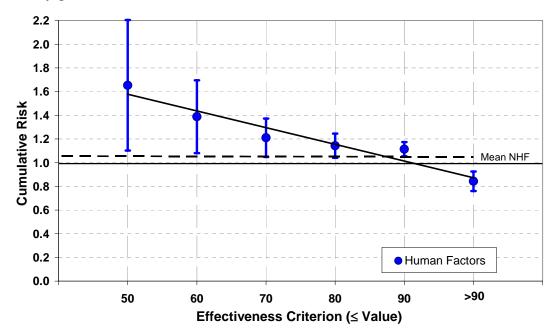


Figure 8. Human Factors Accident Risk by Criterion Levels of Effectiveness Aggregated for Five Railroads

Table 1 summarizes the results in Figure 8, along with the percent of work time and accidents at each criterion level. Several things are clear. First, no specific threshold defines the beginning of the relationship between estimated effectiveness and accident risk. At any criterion level of effectiveness below a score of 90, a significant, albeit a small, increase in risk occurs. Effectiveness above a score of 90 is actually protective; accident likelihood is significantly less than chance for both human factors and nonhuman factors accidents. An effectiveness level of 90 is meaningful; the average daytime worker who consistently gets 8 hr of sleep each night will never have an effectiveness score below 90 according to the model used here. If a person works at night or consistently misses some sleep, effectiveness drops below a score of 90, at least some of the time. Here the analysis indicates that whenever workers are below an

effectiveness level of 90, they have a statistically significant increase in accident risk. However, an increase in risk, while detectable, may not be operationally meaningful. For example, Figure 2 indicates that a considerable amount of work occurs at night, which is an operational necessity. Undoubtedly, effectiveness is below a score of 90 for some or all of that time for those workers; however, that degradation is biologically determined. An 11 percent increase in accident risk may be an operational cost inherent to any industry that operates at night. What is more important to note is that fatigue risk gradually escalates with progressively lower levels of effectiveness.

Table 1. Human Factors Accident Cumulative Risk at Various Criterion Levels of Effectiveness

Criterion Effectiveness <sup>#</sup> Score	Human Factors Accident Risk (%)	Percent of Work Time	Human Factors Cases Number (Percent) <sup>+</sup>
> 90	- 16 *	42	259 (35%)
≤ 90	+ 11 *	58	472 (65%)
≤ 80	+ 14 *	35	289 (40%)
≤ 70	+ 21 *	19	166 (23%)
≤ 60	+ 39 *	7	71 (10%)
≤ 50	+ 65 *	2.7	33 (4.5%)

<sup>\*</sup> Significantly different from chance (p < 0.05).

Choosing a criterion level for effectiveness, below which risk is intolerable, is an economic, operational, and safety decision beyond the scope of this report. From a purely statistical point of view, however, a criterion set at 70 and below indicates the point at which human factors fatigue-associated accident risk is reliably greater than chance and the risk for nonhuman factors accidents.

#### 3.4 Other Evidence of Fatigue

The validation and calibration analyses documented in Figures 6 and 8 highlight human factors accidents that occur below an effectiveness score of 70 as the accidents most strongly related to the influence of fatigue or low effectiveness. The relationship between estimated effectiveness and accident risk is merely a correlation. This relationship could be due to fatigue, or it could be due to some other extraneous factor that is associated

<sup>&</sup>lt;sup>#</sup> Effectiveness at accident time based on 30-day work histories processed using the SAFTE biomathematical fatigue model.

<sup>&</sup>lt;sup>+</sup> Human factors cases (two crewmembers per accident) in  $2\frac{1}{2}$  yrs, excluding accidents involving consistent night workers. The percentages above and below 90 sum to 100 percent. The percentages below 90 are cumulative and do not sum to 100 percent.

with this method for determining effectiveness. One way to more strongly implicate fatigue as an important part of this relationship is to identify evidence that the accidents were caused by errors that are typical of a fatigued state. To pursue this evidence, the analysis determined which cause codes were involved in the human factors accidents with effectiveness scores at or below 70 to see if they were the type of accident expected to be related to fatigue. If those accidents were caused by human errors that might reasonably be expected to result from fatigue, then that would constitute further validation of this method and increase confidence that the reliable statistical relationship is not a coincidental result of some extraneous factor(s).

**Table 2. Human Factors Accident Cause Codes Related to Effectiveness ≤ 70** 

Rank	Cause Code	Frequency w/ Effectiveness ≤ 70	Frequency Total	Relative Likelihood	Category Description
1	H400	15	40	1.53	Main Track Authority
2	H700	42	149	1.15	Use of Switches
	All HF	215	880		All Human Factors Cause Codes, Primary and Secondary Total
Rank	Cause Code	Frequency w/ Effectiveness ≤ 70	Frequency Total	Relative Likelihood	Individual Cause Code Description
1	H701	4	4	4.09	Spring switch not cleared before reversing
2	H311	5	7	2.92	Moving cars while loading ramp, hose, or chute not in proper place
3	H019	8	15	2.18	Failure to release hand brake on cars
4	H222	4	8	2.05	Automatic block or interlocking signal displaying other than stop indication–failure to comply
5	H404	7	17	1.69	Train order, etc., failure to comply
6	H504	4	10	1.64	Buffing or slack action excessive, train makeup
7	H221	12	31	1.58	Automatic block or interlocking signal displaying stop indication—failure to comply
8	H605	18	53	1.39	Failure to comply with restricted speed in connection with the restrictive indication of a block or interlocking signal
9	H399	4	12	1.36	Other general switching rules
10	H503	20	60	1.36	Buffing or slack action excessive, train handling

Top 10 most overrepresented  $\leq$  70, excluding those with < 3 occurrences.

Note: A human factors accident could have two human factors cause codes; hence, more cases of human factors cause codes than human factors accidents occurred.

About 400 human factors accidents occurred, and usually two crew members were associated with each. Hence, about 800 cases occurred in which an effectiveness score was associated with a human factors accident, and 880 human factors cause codes were assigned since an accident could have more than one human factors cause code. Table 2 displays those categories of cause codes and individual cause codes that were overrepresented in human factors accidents with an effectiveness score below 70: the proportion of these accidents at an effectiveness score below 70 was greater than their overall proportion for all human factors accidents. The top two rows of Table 2 indicate the only two categories of causes that were overrepresented as a class: main track authority and use of switches. The lower part of the table lists the 10 most overrepresented individual codes, excluding those that were rare (less than three occurrences in the entire sample). Most of these cause codes are associated with human errors that would be expected to increase with fatigue. The analysis cannot determine if these elevated risks are statistically reliable because the sample of events is too small, but this analysis suggests that the relationship between effectiveness and human factors accident risk is meaningful and not due to extraneous factors.

#### 3.5 Interpretation of Effectiveness Scores

The results of this analysis of freight railroad accident risk indicate a reliable relationship between reduced effectiveness and an increased risk of human factors accidents. Below an effectiveness level of 70, the risk of human factors accidents is increased by about 20 percent; below an effectiveness level of 50, it is elevated by 65 percent. In this section, the effectiveness metric is related to various other metrics and sleep histories to provide a context for understanding and appreciating the kinds of circumstances that can lead to reduced levels of effectiveness of this magnitude.

The effectiveness values shown in Table 3 were derived from the SAFTE model and are based on an average person getting 8 hrs of sleep, awakening at 0700 hr, and remaining awake for the amount of time specified. A lapse is defined as an excessively long reaction time caused by loss of alertness or a micro-sleep. Lapse likelihood is the ratio of the expected frequency of lapses at an effectiveness level to the frequency of lapses of a well-rested person during a normal work day. A lapse likelihood score of 8 means that expected lapses are 8 times more frequent than for a well-rested person (Hursh, et al., 2004). The effects of wakefulness are based on studies of laboratory subjects who were kept awake after a full night of sleep and tested repeatedly on cognitive tests throughout the period of wakefulness (Angus and Heslegrave, 1985; Belenky, et al., 1994). The blood alcohol concentration (BAC) equivalence is based on studies comparing the effects of alcohol and sleep deprivation on performance on a driving simulator (Arnedt, Wilde, Munt, and MacLean, 2001) and cognitive test performance (Dawson and Reid, 1997). As an example, Table 3 shows that at an effectiveness score of 70, lapses are 5 times more likely than for a well-rested person and that this score is the equivalent of being awake for 21 hr after awakening at 0700 hr or having a BAC of 0.08.

Table 3. The Relationship Among Various Effectiveness Scores and Other Meaningful Metrics: Likelihood of a Lapse, Continuous Hours Awake, and BAC

Effectiveness Score	Lapse Likelihood	Hours Awake (Hr:Min)	BAC Equivalent
98	0.2	14:00	
94	1.0	15:10	
90	1.5	16:00	
80	3	18:00	
77	4	18:30	0.05
70	5	21:00	0.08
69	5.4	22:00	
60	8	40:50	
50	12	42:30	
40	18	64:00	

Another way to understand effectiveness values is to consider the kinds of sleep patterns and times of day that can lead to different levels of effectiveness. The values in Table 4 were derived from the SAFTE model and are based on an average person awakening at 0700 hr after getting the hours of sleep shown in the first column or losing the amount of sleep (relative to 8 hr per day) in the second column. Performance at 1600 hr reflects near optimal performance for a person on that schedule; performance at 0400 hr reflects the combined effects of prior sleep, time of day, and 21 hr of wakefulness. The values at 0400 hr do not consider benefits of an evening nap; for example, a 2-hr nap at 2000 hr improves performance at 0400 hr by 4-6 percent.

Table 4. The Effects of Various Daily Sleep Patterns on Effectiveness Estimates at 1600 hr and 0400 hr. Three Schedules: 1, 2, 7 Days at the Specified Sleep Level

Prior Daily	Prior Daily	Effectiveness Score After:					
Sleep	Sleep Loss (Hr)	One Day		Two Days		Seven Days	
(Hr)		1600 hr	0400 hr	1600 hr	0400 hr	1600 hr	0400 hr
8	0	97	70	97	70	97	70
7	1	96	69	95	68	93	67
6	2	94	68	92	66	88	63
5	3	92	65	89	62	82	57
4	4	90	63	84	58	72	48
3	5	87	59	78	51	57	34
2	6	83	55	70	42	*	*
1	7	78	49	58	30	*	*
0	8	73	43	46	15	*	*
* No data available for these conditions							

Relating Table 4 to Table 1, the analysis shows that if a person gets less than 8 hr sleep on a regular basis, then effectiveness at 0400 hr (the circadian minimum) will be below a score of 70, and accident risk will be elevated by at least 21 percent. If the person gets

less than 4 hr sleep, then effectiveness at 0400 hr is below a score of 65 in 1 day, less than a score of 60 in 2 days, and less than a score of 50 in 7 days, at which point accident risk is elevated by 65 percent. After 7 days of 4 hr sleep per day, effectiveness at the circadian peak (about 1600 hr) is nearing a score of 70 or an elevated risk of 14-21 percent in the day time.

#### 3.6 Interpretations and Limitations

This study was designed to test the validity of biomathematical models of fatigue, here the SAFTE model and the FAST software, as tools for evaluating schedules for increased risk of fatigue-related accidents. The data considered–accidents and the 30-day work histories that preceded them—were not a random sample of all locomotive crews in freight rail operations. Hence, the levels of effectiveness calculated by the model should not be interpreted as representative of the railroad work force in general. The study was not designed to determine the extent of fatigue in the freight rail industry. The objective was to test the validity of the mathematical model and to calibrate the model so that one could relate the effectiveness estimates to an estimate of elevated accident risk.

While the biomathematical model used is typical of other sleep and performance models and is designed to simulate the effects of circadian rhythms and sleep patterns on fatigue, it is important to emphasize that the data that drove these predictions were entirely related to work schedule and opportunities to obtain sleep. The study did not measure sleep duration and quality directly, nor did the study directly measure cognitive performance. Operationally, this study related work schedule to accident risk using a sophisticated mathematical algorithm, but it did not directly measure fatigue or the performance manifestations of fatigue. The relationships observed are entirely correlational, not causal. It is theoretically conceivable that some other factor correlated with the same work schedule factors considered by the fatigue model was the real operational cause of the elevated accident risk. That possibility cannot be entirely eliminated with the sorts of data used in this study. In practice, however, that logical possibility may not matter if elimination of those work schedule risk factors results in improved safety. The encouraging results of this study provide justification for such an experiment.

The inferences of sleep opportunity are based on a sleep estimation algorithm called AutoSleep, and the calculated effectiveness values were partly dependent on the settings of AutoSleep. Studies are currently underway to measure actual sleep under typical railroad schedules, and those findings may lead to improvements in the sleep estimation accuracy of AutoSleep, which may lead to greater accuracy in estimates of the true risk of fatigue.

The accident risk levels that were found must be carefully limited to the operational factors that were considered: work schedule, time of day, call times, and commute times. Accidents that were estimated to have occurred with relatively high levels of effectiveness could have been fatigue related if the fatigue were caused by factors not considered by the analysis, such as sleep disorders, poor sleep hygiene, unusual circumstances in the operator's life, such as illness or family demands, or unusual sleep needs. This sort of error might be called a miss, a case of fatigue that was misidentified as not fatigue.

Even in cases where the fatigue model judged that fatigue could have occurred, it is possible that the human factors accident was caused by some other human error, such as distraction or poor training. This sort of error might be termed a false alarm, a case in which the analysis suspected fatigue as the cause but was really some other cause. These two sorts of errors, misses and false alarms, are typical of discriminations involving noise, unidentified factors that lead to error in judgment. When making such judgments, the user of the model must weigh the relative costs associated with the two kinds of errors and set a criterion value for using the model, based on Table 1, that is operationally, economically, and socially optimal. Making that decision is beyond the scope of this study.

Furthermore, fatigue models based entirely on the sorts of data used in this study are not suitable for determining the fatigue of particular individuals because of the kinds of inaccuracies discussed above, especially the wide variation in individual sleep requirements and absence of specific information on individual sleep habits, health, and circumstances. On the other hand, by adding such additional information to a fatigue analysis, it is possible to reduce the error sufficiently to use fatigue models as tools in accident investigations, combined with other important information about the performance of the individuals. That sort of detailed analysis, however, was clearly beyond the scope of this study.

Finally, it is important to properly interpret the risk values in Table 1. Those values represent the percent change in human factors accidents relative to the expected distribution of accidents based on work time effectiveness. The results indicate that the lower the effectiveness level, the greater the elevation of human factors accident risk. The percent changes in risk (second column, Table 1) are not to be confused with the percentages of human factors accidents associated with that level of risk (fourth column, Table 1). The percentages in second column represent the increase (or decrease) in the risk of fatigue-associated human factors accidents compared to chance when at a particular level of effectiveness.

The overall probability of a particular railroad worker on a particular work shift having a human factors accident is extremely low. This study only considered work histories related to accidents, not all work histories. Consequently, an overall probability of fatigue-caused human factors accidents in freight service cannot be calculated, but that probability will be a number with four or five zeros after the decimal point.

## 4. Summary and Conclusions

#### 4.1 Summary

This report summarized the results of a project sponsored by FRA's Office of Research and Development and Office of Safety to develop a methodology to validate and calibrate biomathematical fatigue models for use as fatigue management tools. This study assessed whether a fatigue model can predict an increased risk of human factors caused accidents. Using 2½ years of accident data from five Class I freight railroads, the fatigue model used 30-day work histories before a sample of human factors and nonhuman factors accidents to determine the relationship between accident risk and crew effectiveness (the inverse of fatigue). This report summarized the scope of that effort, the methods used to conduct the analysis, and the results of the analysis to date. The study considered approximately 2800 crew member work histories associated with 1400 accidents, 400 of which involved human factors errors. A biomathematical fatigue model (SAFTE) evaluated a total of over 1 million 30-minute work intervals, covering over 57,000 work starts, for predicted effectiveness. The 2800 work histories served as a basis for determining the exposure level to various levels of fatigue in these locomotive crews. In addition, the fatigue model calculated estimated effectiveness at the time of each accident, again based on 2800 estimates of effectiveness at the time of the accidents. The following summarizes the results of the study.

Human factors accidents follow a circadian pattern that is significantly correlated with the circadian rhythm of a fatigue model (r = 0.71, p < 0.05). The same model rhythm is not correlated with nonhuman factors accidents. The maximum human factors accident risk due to time of day alone was increased less than 20 percent (Figure 3), while the maximum increase in accident risk due to reduced effectiveness (fatigue) was 65 percent (Figure 6), reflecting the combined effects of time of day and sleep opportunities.

- Validation. The fatigue model met two validation criteria.
  - 1. Accident risk was significantly correlated with effectiveness for human factors accidents (r = -0.93, p < 0.01, Figure 6).
  - 2. The model distinguished between human factors and nonhuman factors accidents. Nonhuman factors accident risk was not correlated with estimated effectiveness (r = -0.14, p > 0.05, Figure 7), and the relative risk of human factors and nonhuman factors accidents were significantly different at low effectiveness levels ( $\chi^2 = 7.201$ , p < 0.01).
- Calibration. Effectiveness has a reliable correspondence to increases in accident risk.
  - 1. Risk is a relatively smooth increasing function with decreases in effectiveness
  - 2. Above an effectiveness score of 90, risk is significantly reduced relative to chance.
  - 3. Below an effectiveness score of 70, risk is significantly elevated relative to nonhuman factors risk.

- 4. Overall, railroad workers in this study spent about 40 percent of work time above an effectiveness score of 90 and about 20 percent of work time below an effectiveness score of 70.
- An effectiveness score of 90 is the minimum level for a person getting 8 hrs sleep per day and awake from about 0700 to 2300 hr (day shift).
- An effectiveness score of 70 is the minimum level for a rested person after being awake for 21 hr at 0400 hr.
- An effectiveness score of 70 is about equal to the effects of 0.08 BAC and lapse likelihood five times greater than a well-rested person during the daytime.

An analysis of the cause codes associated with the accidents that occurred at effectiveness at or below an effectiveness score of 70 indicated an overrepresentation of human factors errors associated with main track authority and use of switches. Most of the overrepresented individual cause codes reflect the kinds of operator errors one might expect of persons who are fatigued. This finding confirms that the relationship between effectiveness and human factors accident risk is meaningful and not a circumstantial coincidence.

#### 4.2 Conclusion

This project established that a biomathematical fatigue model can be used to assess how much work schedule factors can contribute to increased fatigue and an elevated risk of railroad accidents. The virtue of having a validated fatigue model, especially if it is calibrated to accident likelihood, is that a carrier could use it to do a self-assessment of fatigue across its system. By evaluating work histories on a terminal by terminal basis and using the scores from the model as a metric, the carrier could determine which terminals are experiencing schedules that might be generating increased risk of fatigue in train operators. Perhaps none of the terminals have a problem, perhaps just a few. In any case, the carrier would be in a position to use this objective assessment as a way to focus its fatigue management efforts where the greatest payoff would be expected. This study provides evidence that such a strategy, using a validated fatigue model, can identify work schedule-related fatigue factors that contribute to an elevated risk of accidents.

#### 5. References

- Abelson, R.P. (1995). Statistics as principled argument. Hillsdale, NJ: Lawrence Erlbaum.
- Angus, R., and Heslegrave, R. (1985). Effects of sleep loss on sustained cognitive performance during a command and control simulation. *Behavior Research Methods, Instruments, and Computers*, 17, 1, 55-67.
- Armitage, P., and Berry, G. (1994). *Statistical methods in medical research* (3<sup>rd</sup> edition). London: Blackwell Science.
- Arnedt, J.T., Wilde, G.J., Munt, P.W., and MacLean, A.W. (2001). How do prolonged wakefulness and alcohol compare in the decrements they produce on a simulated driving task? *Accident Analysis and Prevention*, 33, 3, 337-44.
- Belenky, G., Penetar, D., Thorne, D., Popp, K., Leu, J., Thomas, M., Sing, H., Balkin, T., Wesensten, N., and Redmond, D. (1994). The effects of sleep deprivation on performance during continuous combat operations. In B.M. Marriott (Ed.), *Food components to enhance performance* (pp.127-135). Washington, DC: National Academy Press.
- Belenky, G., Wesensten, N.J., Thorne, D.R., Thomas, M.L., Sing, H.C., Redmond, D.P., Russo, M.B., and Balkin, T.J. (2003). Patterns of performance degradation and restoration during sleep restriction and subsequent recovery: A sleep dose-response study. *Journal of Sleep Research*, 12, 1, 1-12.
- Conover, W.J. (1999). *Practical nonparametric statistics* (3<sup>rd</sup> edition). New York: John Wiley and Sons.
- Dawson, D., and Reid, K. (1997). "Fatigue, alcohol and performance impairment." *Nature 388*, 23.
- Hursh S.R., Balkin, T.J., Miller, J.C., and Eddy, D.R. (2004). The fatigue avoidance scheduling tool: Modeling to minimize the effects of fatigue on cognitive performance. *SAE Transactions*, 113, 1, 111-119.
- Hursh S.R., Redmond, D.P., Johnson, M.L., Thorne, D.R., Belenky, G., Balkin, T.J, Storm, W.F., Miller, J.C., and Eddy, D.R. (2004). Fatigue models for applied research in warfighting. *Aviation, Space and Environmental Medicine*, *75*, 3, Suppl.: A44-53.
- Lowry, R. (1998-2006). VassarStats: Web site for Statistical Computation. http://faculty.vassar.edu/lowry/VassarStats.html
- Newcombe, R.G. (1998). Interval estimation for the difference between independent proportions: Comparison of eleven methods. *Statistics in Medicine*, *17*, 873-890.
- Pollard, J. K. (1996). *Locomotive engineer's activity diary* (Report No. DOT/FRA/RRP-96/02). Washington, DC: U.S. Department of Transportation.
- Van Dongen, H.P.A. (2004). Comparison of mathematical model predictions to experimental data of fatigue and performance. *Aviation, Space and Environmental Medicine*; 75, 3, Suppl.: A15-36.
- Wilson, E.B. (1927). Probable inference, the law of succession, and statistical inference. *Journal of the American Statistical Association*, 22, 209-212.

## **Abbreviations/Acronyms**

AAR Association of American Railroads

BAC blood alcohol concentration

BNSF Railway

CSX CSX Transportation Inc.

FAST Fatigue Avoidance Scheduling Tool

FRA Federal Railroad Administration

KCS Kansas City Southern Railway Company

NS Norfolk Southern Railway Company

SAFTE Sleep, Activity, Fatigue, and Task Effectiveness Model

UP Union Pacific Railroad